



Summary

We provide an example modification to traditional reinforcement learning we call observational dropout, whereby we artificially constrain the probability that an agent is allowed to observe its real environment at each step. In doing so, we can coerce an agent into learning a world model to fill in the observation gaps seen during reinforcement learning without having to explicitly train the world model via teacher forcing. In this way, the policy and world model can be trained purely from the task reward. Visit <u>https://learningtopredict.github.io/</u> for videos and code.

Warmup - Balance Cartpole

Are there any situations in which a "world model" can be learned with total observational dropout?

$$\mathcal{L} = \frac{1}{2}(M+m)\dot{x}^2 + \frac{1}{2}mL^2\dot{\theta}^2 - mL\cos(\theta)\dot{\theta}\dot{x} - mgLc$$
$$\begin{pmatrix} \dot{\theta} \\ \ddot{\theta} \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ \frac{g}{L} + \frac{u_1}{ML} & \frac{u_2}{ML} \end{pmatrix} \begin{pmatrix} \theta \\ \dot{\theta} \end{pmatrix} \sim \begin{pmatrix} a & b \\ c+u_1 & d+u_2 \end{pmatrix}$$

- 1. A world model can be learned that produces a valid policy without needing a forward predictive loss
- 2. A world model need not itself be forward predictive (at all) to facilitate finding a valid policy, and
- 3. The inductive bias intrinsic to one's world model almost entirely controls the ease of optimization of the final policy.

Observational Dropout

At every timestep, with probability p the agent sees the ground truth environment, and with probability (1-p) sees a representation of its environment as determined by its world model.

For this work, world models are constructed to have the same input-output dimensionality as the base environment, and are simple MLP networks.

Learning to Predict Without Looking Ahead: **World Models Without Forward Prediction** C. Daniel Freeman¹ Luke Metz¹ David Ha¹

cdfreeman@google.com

Swingup Cartpole

cos(heta)



Task rewards for policies trained jointly with world model (left) and trained entirely within a learned world model (right) as a function of peek p.



Task reward for policies trained with fully connected and convolutional world models. Inductively biased world model architecture dominates.



Task reward for policies trained jointly with VAE-based world model (left), and using the representation learned by a world model (right).

¹Google Brain

Deploying Policy Learned in World Model to Actual Environment





High level points:

- using only a task reward.
- be useful.

Open problems: • A better way to formalize "minimal world model" Discrete environments are hard

- Optimal p?





• Observational dropout provides a way to train world models

World models need not be (nor should be?) "pixel perfect" to